



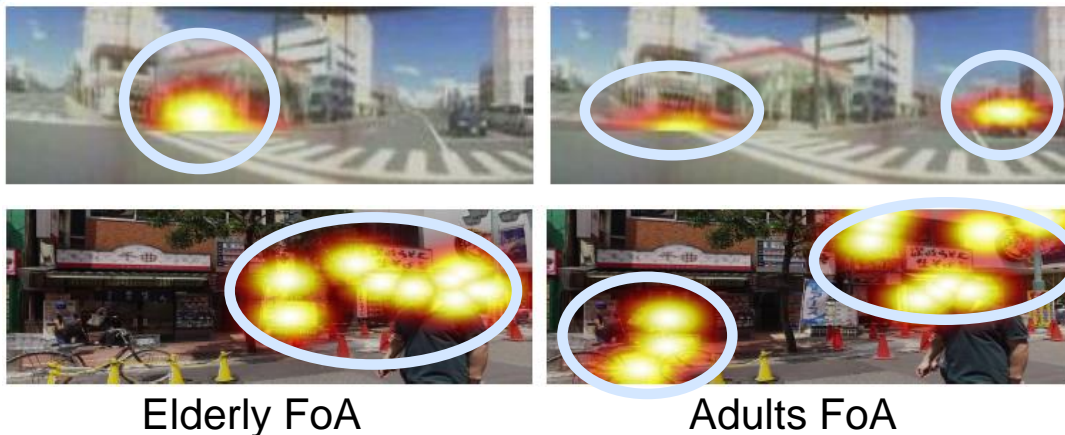
# Translating Adult's Focus of Attention to Elderly's

Onkar Krishna   Go Irie   Takahito Kawanishi   Kunio Kashino   Kiyoharu Aizawa

NTT Corporation   The University of Tokyo

# Motivation

- FoA is a region of an image/video which attracts our attention
- FoA of elderly is significantly different from other age-groups **due to aging!**
- Existing FoA prediction models **fail to predict elderly FoA** as they are developed and evaluated on adults' eye-gaze data.



**Our goal is to propose an approach to predicting the elderly FoA**  
for assisting their daily activities, such as driving, walking, and searching.

# Challenges and Approach

- **Straightforward Approach** – training an FoA prediction model on elderly's data.
  - ☹ Collecting a sufficient amount of data from elderly's is more challenging than adult, due to their physical or health conditions.
- **Assumption:** Correlation between adult's and elderly's FoAs can be characterized by the scene they are viewing.
- **Our Proposal:** Image-to-image translation from adult's FoA to elderly's.
  - 😊 Leveraging well-trained models for adult's FoA for data efficient training.



# Problem Setting

- **Input:** Sequence of video frames
- **Output:** Eldery's FoA maps for input videos viewed in two different scenarios



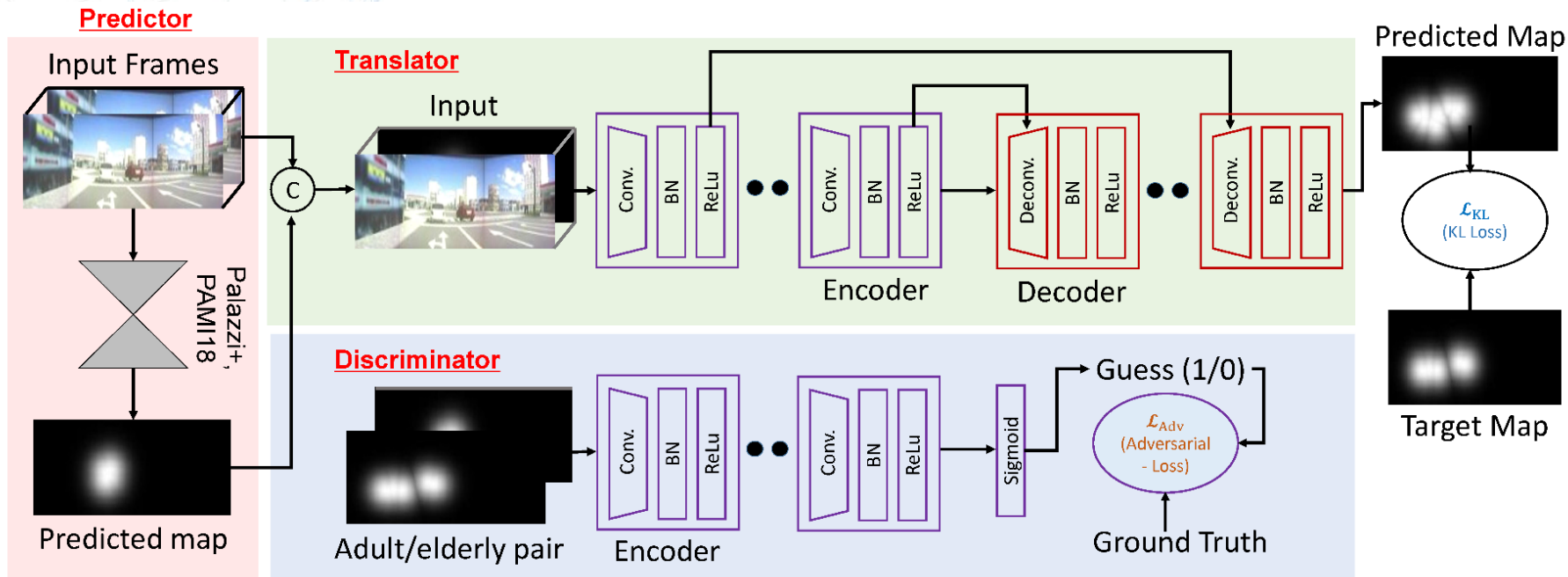
Driving  
car driving scenario



Street Video  
free viewing scenario during street walking

# Network Architecture

Our model has **predictor network** and **translator network** trained w/ support of auxiliary **discriminator network**



# Loss Function Design

Our model is trained by minimizing both **Kullback-Leibler divergence** and **Adversarial loss**

$$\min_D \max_T \mathcal{L}_{\text{Adv}}(T, D) - \gamma \mathcal{L}_{\text{KL}}(T)$$

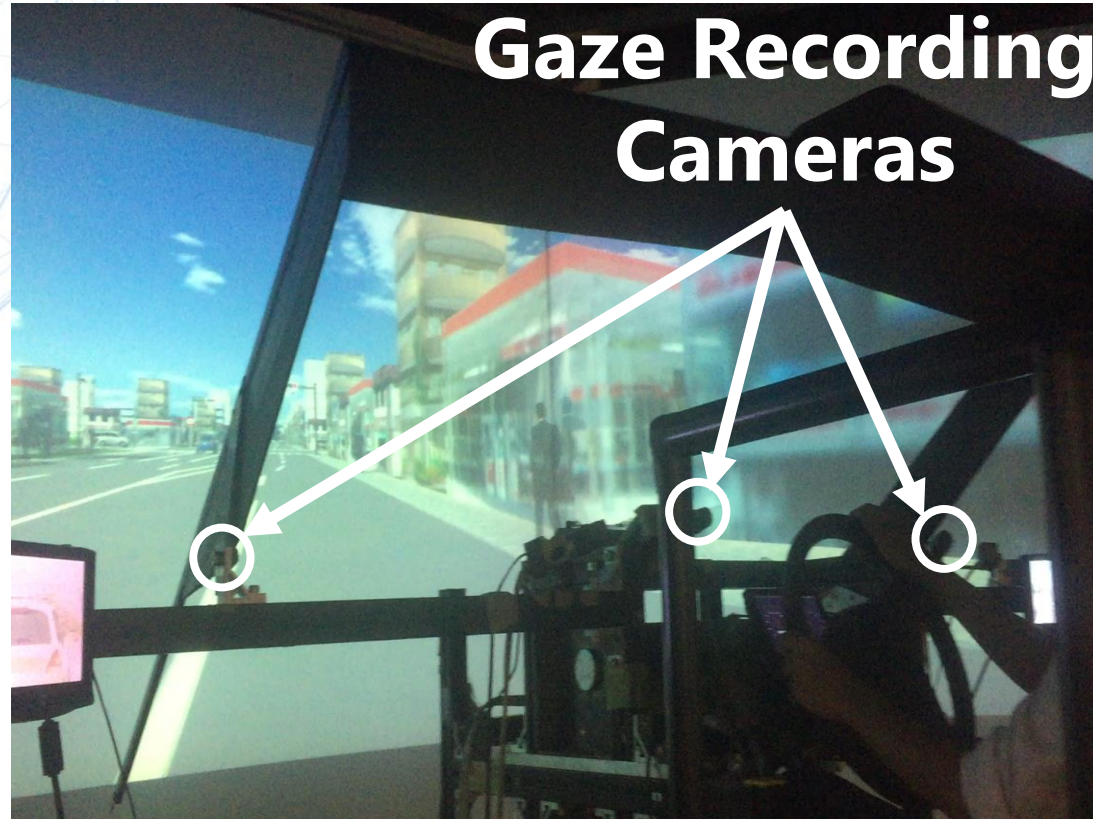
$\mathcal{L}_{\text{KL}}$  loss requires the translator network  $T$  to output the ground truth FoA

$$\mathcal{L}_{\text{KL}}(T) = \sum_n \sum_i e_{n,i}^* (\log(e_{n,i}^*) - \log(e_{n,i}))$$

$\mathcal{L}_{\text{Adv}}$  loss approximate the joint probability distribution of adult and elderly FoA

- We construct **Driving dataset** and **Street Video dataset** covering both adults and elderly.
- **Participants**
  - 18 participants (adults and elderly)
  - Adult's mean age **26 years**, elderly mean age **75 years**.
- **Eye-gaze for Driving**
  - Collected fixations of each participant while driving on a car simulator in real time
  - 9,713 FoA maps correspond to the 9,713 frames (train-test split 7,716/1,997)
- **Eye-gaze for Street Video**
  - Collected fixations of each participant while watching street-walking videos displayed on a monitor
  - 4,425 FoA maps correspond to the 4,425 frames (train-test split 3,532/893)







# Result on Driving Dataset

Algorithm	CC $\uparrow$	SIM $\uparrow$	KL-div. $\downarrow$	Time (sec.) $\downarrow$
[Wang+, 19]	0.13	0.22	5.60	6.31
[Wang+, 15]	0.09	0.26	4.90	6.43
[Cornia+, 16]	0.26	0.42	9.97	2.71
[Palazzi+, 18]	0.64	0.53	4.06	7.48
[Palazzi+, 18] (fine tuned)	0.66	0.55	3.89	7.48
<b>Ours</b>	<b>0.91</b>	<b>0.79</b>	<b>0.80</b>	7.56

**Ours achieves huge performance gain with this slight expense of run time compared to the base method**

# Result on Street Video Dataset

Algorithm	CC $\uparrow$	SIM $\uparrow$	KL-div. $\downarrow$	Time (sec.) $\downarrow$
[Harel+, 07]	0.24	0.49	0.82	4.10
[Itti+, 2000]	0.22	0.47	1.00	6.33
[Jiang+, 18]	0.27	0.46	2.21	9.23
[Cornia+, 18]	0.27	0.47	1.36	<b>2.74</b>
[Cornia+, 18] (fine tuned)	0.58	0.57	<b>0.58</b>	<b>2.74</b>
<b>Ours</b>	<b>0.72</b>	<b>0.71</b>	0.94	2.93

**Ours outperform all the baselines for CC and SIM scores with slight expense of run time**

# Qualitative Results

## Driving Dataset



## Street Video Dataset



**Ours model accurately mimics the gaze of elderly people!**

# Summary

**We introduced a deep image translation framework for predicting the elderly's FoA.**

- *Accuracy*: our model can accurately mimic the elderly FoA while driving and street walking which can be useful in assisting elderly.
- *Novel Training*: adversarial training together with KL-divergence loss allows us to reach state-of-the art performance.